

# Variational Mixture Models with Gamma or inverse-Gamma components

A. Llera<sup>1</sup>, D. Vidaurre<sup>2</sup>, R.H.R. Pruim<sup>1</sup>, C. F. Beckmann<sup>1</sup>.

Technical report

1- Donders Institute for Brain Cognition and Behaviour, Radboud University Nijmegen

2- Oxford centre for Human Brain Activity (OHBA)

## Abstract

Mixture models with Gamma and/or inverse-Gamma distributed mixture components are useful for medical image tissue segmentation or as post-hoc models for regression coefficients obtained from linear regression within a Generalised Linear Modeling framework (GLM), used in this case to separate stochastic (Gaussian) noise from some kind of positive or negative 'activation' (modeled as Gamma or inverse-Gamma distributed). To date, the most common choice in this context it is Gaussian/Gamma mixture models learned through a maximum likelihood (ML) approach; we recently extended such algorithm for mixture models with inverse-Gamma components. Here, we introduce a fully analytical Variational Bayes (VB) learning framework for both Gamma and/or inverse-Gamma components.

We use synthetic and resting state fMRI data to compare the performance of the ML and VB algorithms in terms of area under the curve and computational cost. We observed that the ML Gaussian/Gamma model is very expensive specially when considering high resolution images; furthermore, these solutions are highly variable and they occasionally can overestimate the activations severely. The Bayesian Gauss/Gamma is in general the fastest algorithm but provides too dense solutions. The maximum likelihood Gaussian/inverse-Gamma is also very fast but provides in general very sparse solutions. The variational Gaussian/inverse-Gamma mixture model is the most robust and its cost is acceptable even for high resolution images. Further, the presented methodology represents an essential building block that can be directly used in more complex inference tasks, specially designed to analyse MRI/fMRI data; such models include for example analytical variational mixture models with adaptive spatial regularization or better source models for new spatial blind source separation approaches.

# 1 Introduction

Mixture models are an important and powerful tool in many practical applications thanks to their ability to flexibly model complex data [17]. Mixture models containing Gamma or inverse-Gamma distributed components are interesting due to the positive support of such distributions and are commonly used to provide class-dependent models separating stochastic noise, typically modeled by a close to zero-mean Gaussian distribution, from some kind of activation modeled by a positive support distribution [13, 2, 24]. For example, in medical imaging such models can be used for statistical segmentation of structural images into different tissue types on the basis of measured intensity levels. Also, in functional statistical parametric mapping (where voxels are either activated or not activated), mixture models can be used for post-hoc inference on the regression maps[25].

The most common approach to learn mixture models in general is the expectation maximization EM algorithm (EM) which is used to estimate a maximum likelihood (ML) solution [6]. However, since there is no closed form ML solution neither for the scale parameter of the Gamma nor for the shape parameter of the inverse-Gamma, the problem becomes more complex and typically requires numerical optimization [14, 15, 2, 31]. Numerical optimization must be performed at each iteration of the EM algorithm, making such strategy computationally hard, specially for cases where the number of samples is very high, as e.g. high resolution whole brain MRI data. A common faster alternative uses the method of moment approximation to estimate the parameters of the Gamma or inverse-Gamma components [5, 18, 30]. We denote the algorithm presented in [5] for learning a Gauss/Gamma mixture model as GGM, and the one presented in [18] for learning Gaussian/inverse-Gamma ones as GIM. An alternative to such ML approaches is to consider Bayesian inference. The Bayesian approach provides an elegant way to explore uncertainty in the model and/or to include prior knowledge into the learning process. Furthermore, it provides principled model selection to select the number of components in the mixture model, and it allows to use the learnt components as building blocks of bigger Bayesian inference problems [7]. To the extent of our knowledge, there are sampling algorithms available for the Gamma case [22, 9] and versions providing spatial regularization [28]. However, the sampling strategy can be computationally infeasible for high resolution images, and

specially in cases where the mixture distributions become part of bigger statistical learning problems [21]. Variational Bayes (VB) inference [6] provides instead a more efficient alternative. Although in [29] a variational Gaussian/Gamma mixture model with spatial regularization is presented, the Gamma distribution parameters of the mixture model are learnt using a conjugate gradient numerical optimization procedure. In this work we introduce novel algorithms for learning mixture models with Gamma and/or inverse-Gamma components using an analytic VB approach. While most parameters belong to conjugate distributions and can be estimated easily, learning the shape parameter of the distributions is not so straightforward. For the shape parameters we use unnormalized conjugate priors [10, 19], resorting to Laplace approximations and Taylor expansions to compute the required posterior expectations.

In section 2, we introduce the four considered models and outline the datasets used to evaluate them. In section 2.1, we introduce the basic notation and a brief description of the learning algorithms. Further details are given in the Appendix. In section 2.3, we describe the synthetic data sets we consider for evaluation of the methods. In section 2.4 we describe the resting state fMRI dataset as well as the data processing performed to obtain 4400 spatial maps extracted from 100 subjects rfMRI data. In sections 3.1 and 3.2 we present the results obtained by comparing the two newly proposed algorithms with their maximum likelihood counterparts in both artificial and rfMRI data. Finally, in section 4, we conclude the paper with a brief discussion.

## 2 Methods

We now introduce the methodology and the datasets used to evaluate the different considered models. In section 2.1, we introduced the notation necessary to describe the problem and, in section 2.2, we introduced the two state of the art models alongside their two new Bayesian versions. Then, we introduce the (synthetic and rfMRI) datasets that will later be used to evaluate the four considered models.

## 2.1 The problem

Let  $\mathbf{x} = \{x_1, \dots, x_N\}$ ,  $x_i \in \mathbb{R}$  be an observation data vector. Without algorithmic loss of generality we will reduce derivations to mixture models of three components, so that  $p(\mathbf{x}|\boldsymbol{\pi}, \Theta) = \prod_{n=1}^N \sum_{k=1}^3 \pi_k p_k(x_n|\Theta_k)$ , where  $\Theta = \{\Theta_1, \Theta_2, \Theta_3\}$  are the parameters of the three components and  $\boldsymbol{\pi} = \{\pi_1, \pi_2, \pi_3\}$  are the mixing proportions. One component is used to model stochastic noise which, as usual, is modeled using a Gaussian component:  $p_1(x|\mu_1, \tau_1) = \mathcal{N}(x|\mu_1, \tau_1)$  with  $\mu \approx 0$ . The other two components model independently positive and negative activations. Here we extend the common choice of Gamma distributions to consider also inverse-Gamma components, that means that the positive component  $p_2(x|s_2, r_2)$  can be chosen to be Gamma

$$p_2(x|s_2, r_2) = \mathcal{G}(x|s_2, r_2)$$

or inverse-Gamma distributed

$$p_2(x|s_2, r_2) = \mathcal{IG}(x|s_2, r_2),$$

and the negative component  $p_3(x|s_3, r_3)$  can be negative Gamma

$$p_3(x|s_3, r_3) = \mathcal{G}^-(x|s_3, r_3) = \mathcal{G}(-x|s_3, r_3)$$

or Negative inverse-Gamma distributed

$$p_3(x|s_3, r_3) = \mathcal{IG}^-(x|s_3, r_3) = \mathcal{IG}(-x|s_3, r_3).$$

Regardless of the choice of the distribution,  $s_k$  represents the shape of the distribution; for any Gamma component  $r_k$  denotes the rate parameter, and for the inverse-Gamma ones it denotes the scale parameter. A general graphical representation is presented in the left panel of Figure 1.

## 2.2 The solutions

Learning the model parameters  $\Theta = \{\boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}\}$  is usually achieved through the EM algorithms presented in [5, 18, 30]. These algorithms use the method of moment approximation for the Gamma or inverse-Gamma component parameters

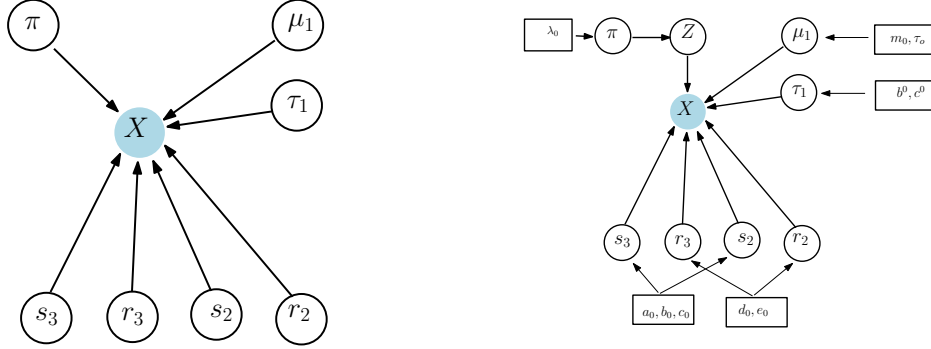


Figure 1: Figure 1 left panel shows a graphical representation of a mixture model with 3 components, one Gaussian and two Gamma and/or inverse-Gamma distributed. The right panel shows such a representation when including prior distributions over the mixture model parameters as well as an indicator variable  $Z$  (see text for more details).

in order to compute the so-called responsibilities and update the expected means and variances analogous to an EM for Gaussian Mixture model [5, 18, 30]. An alternative to maximum likelihood (ML) approaches is to perform Bayesian inference. Defining prior distributions over each parameter, the right panel of Figure 1 shows a graphical representation for such mixture models where the hyper-priors parameters are represented inside the rectangles. We use a Dirichlet prior for the mixing proportions  $\pi$ , a Gaussian prior for the Gaussian mean  $\mu_1$  and a Gamma prior for its precision  $\tau_1$ . For the  $r$  parameters we use a Gamma prior. For the shape parameter of the Gamma we use the unnormalized conjugate prior proposed in [10] and for the inverse-Gamma the prior we recently introduced in [19]. Note that we also introduced an indicator function  $Z$ , so that, for each observation  $x_n$ , we define a latent variable  $\mathbf{z}_n$  as a binary vector with elements  $z_{nk}$ ,  $k \in \{1, 2, 3\}$ , such that  $\sum_{k=1}^3 z_{nk} = 1$  and we define  $\mathbf{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_N\}$ . Consequently, the vector  $\mathbf{z}_n$  has a value of one in the component number to which  $x_n$  belongs. For a given set of initialization parameters and hyper-parameters values, the posterior expectation on the parameters can in most cases be easily computed by evaluating expectations over well-known distributions. However, computing the shapes posterior expectations is not straightforward. Here, we use Laplace approximations and Taylor expansions to approximate the solution. In Appendix A we introduce

for the first time the detailed methodology that allows us to perform VB inference in such models. We will further denote these algorithms as Algorithm 1 or bGGM and Algorithm 2 or bGIM for the Gamma and inverse-Gamma cases, respectively. For ease of notation, we will denote the ML Gaussian/Gamma algorithm presented in [5] as Algorithm 3 or GGM, and the ML Gaussian/inverse-Gamma algorithm presented in [18] as Algorithm 4 or GIM. For completeness, the method of moments identities as well as both ML algorithms are detailed in the Appendices B and C respectively.

### 2.3 Synthetic data

Synthetic dataset I is generated from Gaussian mixture models with three components and different parameter values. One component has always mean zero while the other two have means  $SNR$  and  $-SNR$  respectively, with  $SNR \in \{2, 3, 4, 5\}$ . The variance of all components is always one and we consider three different levels of symmetric activation,  $\pi \in \{[.8, .1, .1], [.9, .05, .05], [.99, .005, .005]\}$ . These three levels of activation will be denoted respectively as sparsity 1, 2 and 3. In general, a stronger activation makes easier the problem; the range of considered activations was chosen to illustrate a range of problems, from easy at 20 % activation to difficult at 1%. The intermediate proportion (10%) is intended to emulate a strong rfMRI activation. At each simulation we generate  $N=10000$  samples (voxels).

We also consider another synthetic dataset, Synthetic dataset II, which is generated similarly to Synthetic dataset I but with mixing proportions  $\pi \in \{[.9, .1, 0], [.95, .05, 0], [.99, .01, 0]\}$ . Thus, Synthetic dataset II contains positive activation but no negative activation.

For each of the synthetic datasets and for each possible of the 12 possible combinations of SNR and mixture proportions, we generated  $N$  samples from such mixture model and we repeated the process 100 times. In all scenarios we fitted mixture models with three components; therefore, Synthetic dataset II is intended to study the performance of the models with a wrong model order.

## 2.4 Resting State fMRI data

We use resting state fMRI (rfMRI) data from 100 healthy controls from the NeuroIMAGE project; this subset of healthy subjects has been previously used in [23]. For specific information on the scanning protocol and parameters of the NeuroIMAGE datasets we refer the reader to [26]. All rfMRI data processing was carried out using tools from the FMRIB Software Library (FSL<sup>1</sup>) [25, 30, 11]. The preprocessing involved removal of the first five volumes to allow for signal equilibration, head movement correction by volume-realignment to the middle volume using MCFLIRT [12], global 4D mean intensity normalization, 6mm full-width-half-maximum (FWHM) spatial smoothing, ICA-AROMA based automatic removal of residual motion artifacts [23], nuisance regression (using mean white matter, CSF time-courses and linear trend as nuisance regressors) and temporal high-pass filtering ( $>0.01$  Hz). For each participant we transformed the rfMRI data to his/her structural image using FLIRT [12], an affine boundary-based registration. Then, we registered the functional data to the 4mm isotropic resolution MNI152 standard space using a non-linear registration procedure (FNIRT [1]).

To delineate a set of group-level spatial components we conducted a temporal concatenated group-ICA on the preprocessed data using MELODIC [4], where the model order was automatically estimated, resulting in a number of 11 components. Individual spatial maps were derived from the group maps using dual regression [3] for a total of  $11 \times 100 = 1100$  spatial maps.

To compare the performance of the models under different image resolutions we also resampled all these spatial maps to 3mm, 2mm and 1mm isotropic resolution MNI152 standard space, using FLIRT [12]. Altogether, we have a total of 4400 spatial maps.

## 3 Numerical Results

In this section we compare the four considered models, bGGM, bGIM, GGM and GIM. In section 3.1, we evaluate the models using the synthetic data reported in section 2.3. In section 3.2, we test them on the statistical maps extracted from resting state fMRI as described in section 2.4. In the remaining, we will denote the

---

<sup>1</sup><http://www.fmrib.ox.ac.uk/fsl>

different models, bGGM, bGIM, GGM and GIM as algorithms 1-4, using the following color code to identify the models: red=bGGM, green=bGIM, pink=GGM and blue=GIM.

### 3.1 Results on synthetic data

The four considered models are evaluated first in terms of the area under the curve (AUC), normalized in the range  $\text{FPR} \in [0, 0.05]$ . In all cases we fitted mixture models with three components. As expected, we observed that all algorithms benefit from a higher SNR and show higher variance at sparser activations (not shown). For each different SNR and mixture proportions, we then compare each pair of models using a paired t-test. In Figure 2 we present histograms reflecting the percentage of times a model was significantly better than any other model (statistical significance is considered for p-values  $< 0.01$ ).

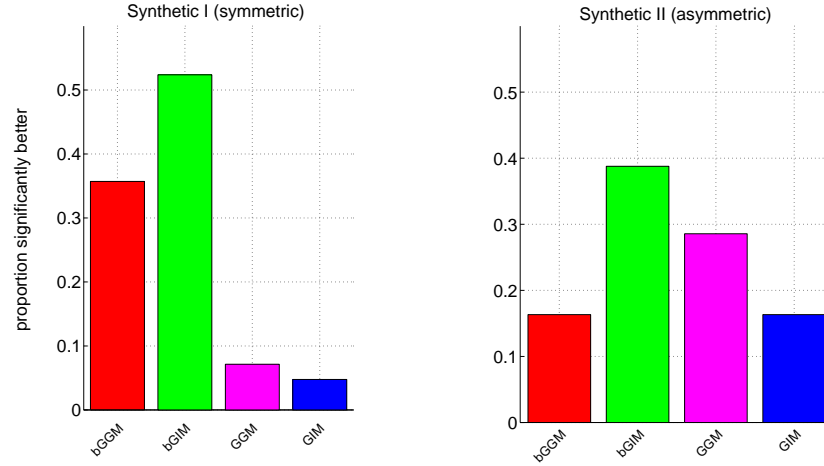


Figure 2: Histogram reflecting the percentage of times each model provides significantly higher normalized AUC than another one. Left pannel shows results in Synthetic dataset 1 (symmetric activation) and the right pannel in synthetic dataset 2 (only positive activation).

The left pannel of figure 2 presents the results obtained on synthetic dataset I (symmetric activation) and the right one on synthetic dataset II (only positive activation). In the case of synthetic dataset I, we observed that VBGGM and VBGIM were the best models. Further, VBGGM was better than VBGIM at the lowest SNR with strongest activations while VBGIM was better in all the other



scenarios. With respect to synthetic dataset II, VBGIM and MLGGM were the best two models and, again, VBGIM was best in most cases with the exception of the low SNR and strong activation cases.

To get a more intuitive idea of the solutions delivered by each model, in Figure 3 we present violin plots of the percentage of positive and negative active voxels provided by each model when considering synthetic dataset I; for visualization the negative proportion is presented as a negative number. The black discontinuous horizontal lines represent the true activation percentages. Most models provide

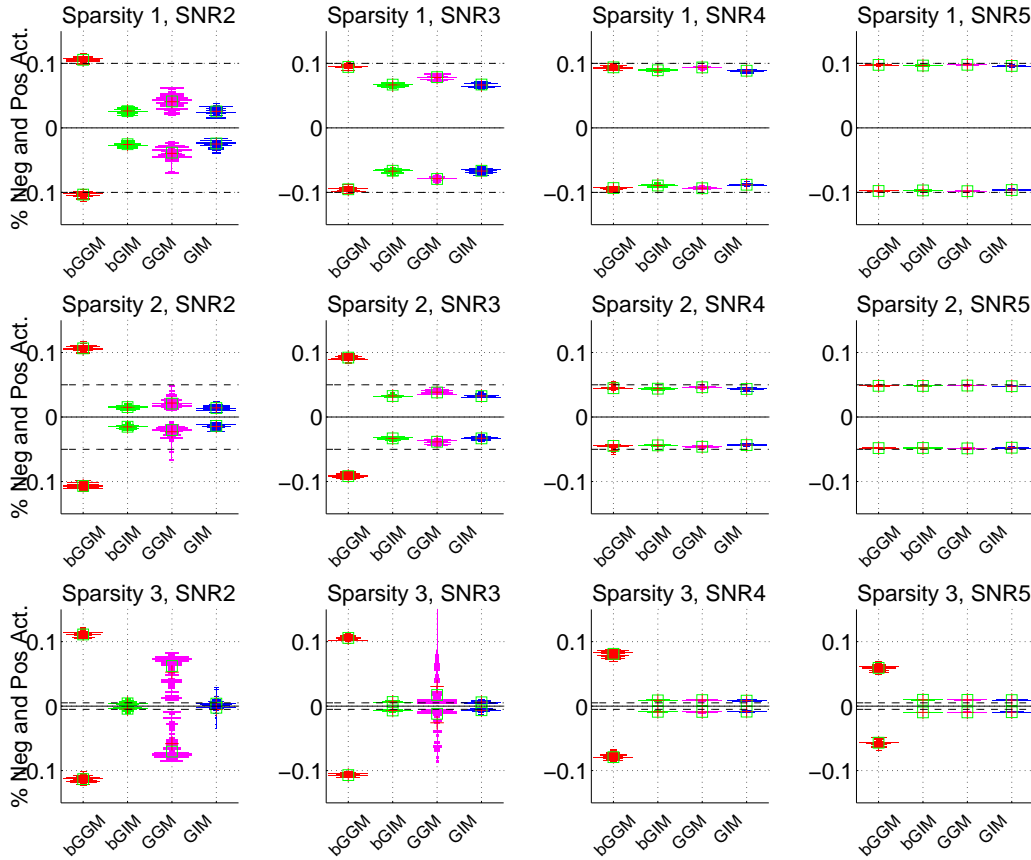


Figure 3: Violin plots of the percentage of positive and negative active voxels of each algorithm (color coded) for Synthetic data I. Each row presents a different symmetric activation levels (or datasets). Each column considers a different SNR. The black discontinuous horizontal line represents the true value.

generally accurate mixing proportions at high SNR. While GIM provides generally very sparse solutions, GGM shows the highest variance in the solutions and overestimates activations specially at low SNR (first column) or sparse activations (last row). Relating the Bayesian models, the bottom row shows that in scenarios where activation is very sparse, the variational Gamma model, bGGM, overestimates activations even at high SNR; the bGIM solution is sparser than the Gamma models and it is very robust as reflected by the low variance in the solutions shown at all SNR and different mixture proportions. Although AUC indicates that bGGM is often an appropriate model, this algorithm overestimates activations at sparser cases. This seemingly contradictory effect occurs because, although the Gamma distribution might overestimate the activation, it still models fairly well the tail of the distribution, which is reflected in the restricted AUC measure. Note that the restricted AUC is a reasonable validation measure when considering fMRI data since more than 5% of false positives would provide meaningless results.

In Figure 4 we present violin plots of the percentage of positive and negative active voxels when considering synthetic dataset II. As before, for ease of visualization, the negative proportion is presented as a negative number in every different scenario. The black discontinuous horizontal line marks the true activations percentage at each dataset. Note that the only difference between this dataset and the previous synthetic data I is that synthetic data II contains no negative activation. Thus, fitting a mixture model with three components to such images could potentially model an unexisting negative activation. Again GIM (blue) provides very sparse positive activation when the activation is strong (first row) but it is also the best estimating the absent activation. On the other hand, bGGM (red) solutions are too dense, modelling non-existing activations specially at low SNR or sparse activations. GGM (pink) shows again the highest variance of all 4 models but it provides a good performance at high SNR even at very sparse activations (bottom rows, right sub-figures). GGM can severely overestimate activations. The bGIM (green) algorithm slightly overestimates the extremely sparse activations (rows 2 and 3) and provides good solutions for the most realistic activation density. Again, bGIM proves to be very robust.

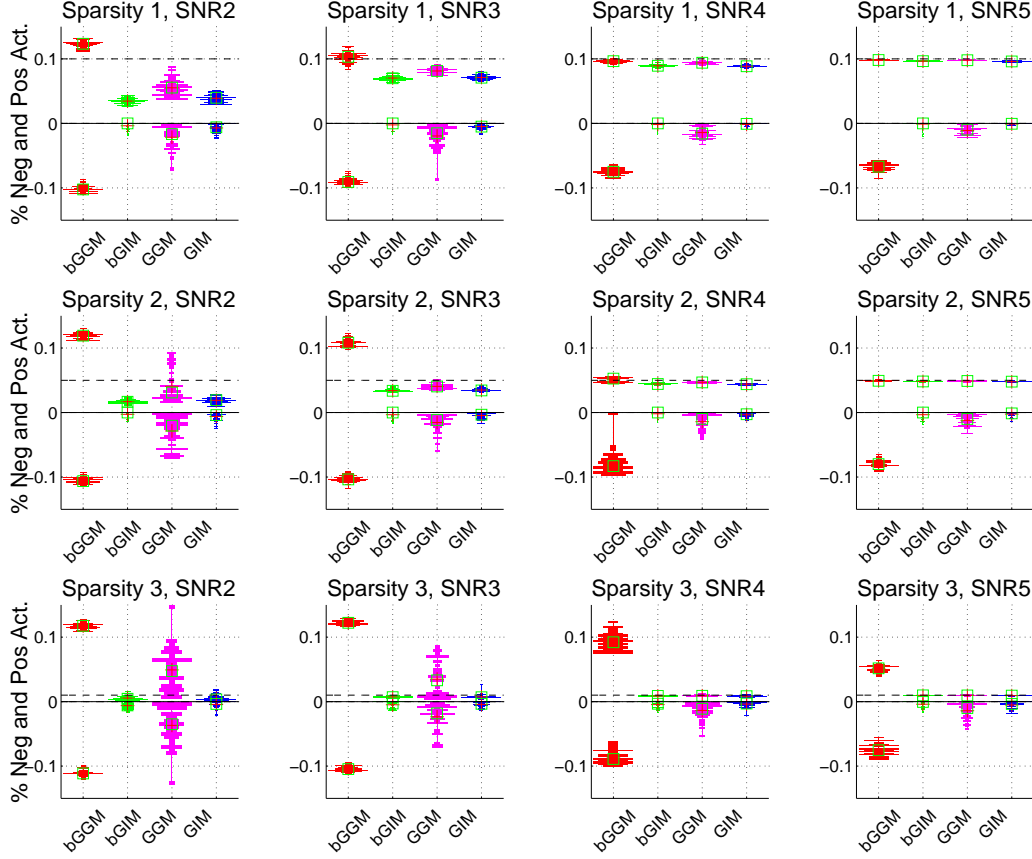


Figure 4: Violin plots of the AUC of each algorithm (color coded) evaluating synthetic data II. Each row corresponds to a different level of sparsity and each column corresponds to a different SNR.

### 3.2 Resting State fMRI data

In this section we compare the four considered models when applying them to the 4400 statistical spatial maps derived from the resting-state fMRI data as described in section 2.4. Each image was masked to remove zero valued voxels and then standardized to zero mean and unit variance. We consider as active those voxels with a probability of activation bigger than a given threshold of  $p=0.5$ . In Figure 5, we present the activation maps obtained by each of the four algorithms when evaluating a pseudo-random spatial map from the 1100 images at 1mm. Color coded are as before. The bGGM model (red) provides the most dense solution,

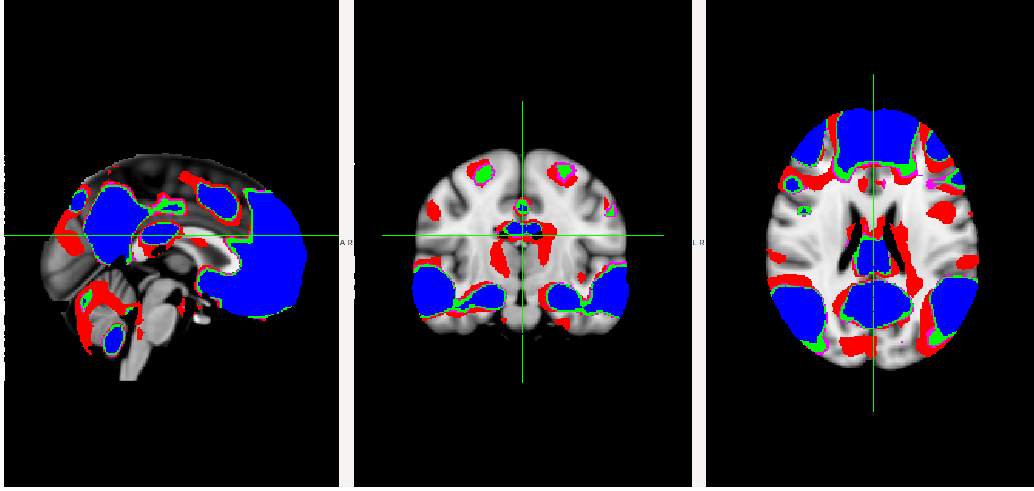


Figure 5: Example of activation maps provided by each algorithm. Color coded are as before.

followed by GGM (pink) and bGIM (green). The sparsest solution is given by GIM (blue). While bGGM provides a much denser solution than the other models, the difference between GGM and bGIM is moderate in this example. The solution provided by GIM is much sparser and it could omit interesting information as can be observed in the middle panel of figure 5; note that the symmetric superior activation reflected by all other models is neglected by GIM.

To summarize the results obtained in the 4400 maps, in Figure 6 we show violin plots on the percentage of active voxels obtained by each algorithm (x-axis and color coded) at four different image resolutions as showed on each subfigure title. The proportion of negative active voxels is presented as a negative number. All models agree in having more positive than negative activation. Independently of the image resolution, the GIM model provides the sparsest images followed by the bGIM; the most dense solutions are given by bGGM. The high variance in the GGM estimations shows that GGM probably overestimated the activation maps.

Another important factor to keep on mind is the computational cost of each algorithms. In figure 7 we show violin plots of the computational cost (in second) taken by each algorithm. From left to right we show the statistics obtained on the 1100 maps obtained at 1, 2 3, 4 mm MNI space respectively. We observe that bGGM is always the fastest followed by GIM and bGIM. The GGM is clearly the most computationally demanding with a cost distribution showing high vari-

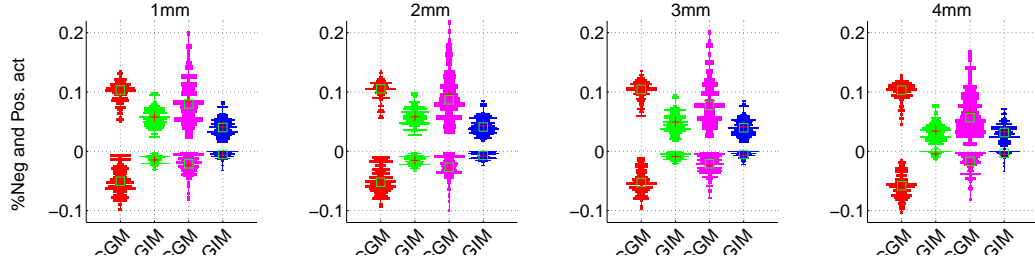


Figure 6: Each subplot shows violin plots of percentage of active voxels. The first row shows results for the positive component and the second row for the negative component.

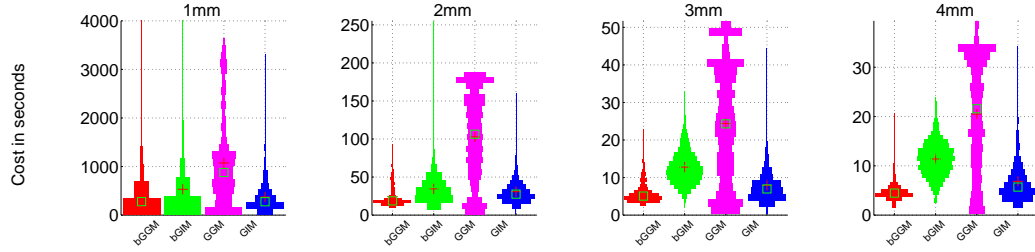


Figure 7: Each subplot shows violin plots of the algorithmic computational costs in seconds for a given image resolution.

ance; the cost is specially large for high image resolutions (left subfigure). The bGIM cost distribution is compact and its cost remains acceptable even for high resolution images.

## 4 Discussion

In this paper we reviewed the state-of-the-art algorithms for learning the parameters of mixture models containing Gamma (GGM) and inverse-Gamma components (GIM), and we introduced novel analytical variational Bayes learning procedures for these mixture models, denoted as bGGM and bGIM respectively. The updates for most model parameters are obtained using standard variational Bayes techniques; for the most involved ones we used Laplace approximations and Taylor expansions to compute the required expectations. We validated the performance of the algorithms in different simulated scenarios and extensive rfMRI data. As is usually done on rfMRI data, we fitted mixture models with three components (for

both real and synthetic data).

We observed that, in general, GIM provides the sparsest solutions, followed by bGIM, GGM; bGGM provides too dense solutions. The GGM solutions showed the highest variance of the four models and overestimated activations with respect to other models in the context of rfMRI data. GIM generally underestimates activations and the bGIM model provides an interesting intermediate solution. Evaluating the models using paired t-tests, bGIM turned out to be the best model in most cases. When considering the computational cost we observed that bGGM is the fastest model closely followed by GIM and bGIM. All models enjoy significant computational advantages with respect to the previous state-of-the-art GGM, the difference becoming dramatic for high image resolutions.

Put together, the bGIM model is an excellent candidate to replace GGM in many neuroimaging tasks. The presented variational methodology also allows the inclusion of Gamma or inverse-Gamma components in more complex inference problems, for example extending VB mixture models for image segmentation [16, 8, 27] to mixtures containing non-Gaussian components. In particular, it can be used to extend the work of [29] to substitute the costly numerical optimization procedure for the Gamma parameters estimation. Another important use of the presented models is in the context of variational ICA decompositions with a Gauss/Gamma or Gauss/inverse-Gamma source model. This assumption on the source model can enhance the sensitivity of the method by placing the source model assumption inside the learning procedure instead of as a post-hoc process.

## References

- [1] J. Andersson, M. Jenkinson, and S. Smith. Non-linear registration, aka spatial normalisation. *Oxford, United Kingdom*, 2007.
- [2] A. Balleri, A. Nehorai, and J. Wang. Maximum likelihood estimation for compound-gaussian clutter with inverse gamma texture. *Aerospace and Electronic Systems, IEEE Transactions on*, 43(2):775–779, April 2007.
- [3] C. Beckmann, C. Mackay, N. Filippini, and S. Smith. Group comparison of resting-state fmri data using multi-subject ica and dual regression. *Neuroimage*, 47, S148, 2009.

- [4] C. Beckmann and S. Smith. Probabilistic independent component analysis for functional magnetic resonance imaging. *IEEE Trans Med Imaging*, 23(2):137–152, 2004.
- [5] C. Beckmann, M. Woolrich, and S. Smith. Gaussian / Gamma mixture modelling of ICA/GLM spatial maps . In *9th Int. Conf. on Functional Mapping of the Human Brain*, 2003.
- [6] C. M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer, 1 edition, 2007.
- [7] R. Choudrey. *Variational Methods for Bayesian Independent Component Analysis*. PhD thesis, University of Oxford.
- [8] C. Constantinopoulos and A. Likas. *Computer Analysis of Images and Patterns: 12th International Conference, CAIP 2007, Vienna, Austria, August 27-29, 2007. Proceedings*, chapter Image Modeling and Segmentation Using Incremental Bayesian Mixture Models, pages 596–603. Springer Berlin Heidelberg, Berlin, Heidelberg, 2007.
- [9] K. Copsey and A. Webb. Bayesian gamma mixture model approach to radar target recognition. *Aerospace and Electronic Systems, IEEE Transactions on*, 39(4):1201–1217, Oct 2003.
- [10] D. Fink. A compendium of conjugate priors, 1997.
- [11] M. Jenkinson, C. Beckmann, T. Behrens, M. Woolrich, and S. Smith. Fsl. *Neuroimage*, 62:78290, 2012.
- [12] M. Jenkinson and S. Smith. A global optimisation method for robust affine registration of brain images. *Med. Image Anal.*, 5:143156, 2001.
- [13] A. Khalili, D. Potter, P. Yan, L. Li, J. Gray, T. Huang, and S. Lin. Gamma-normal-gamma mixture model for detecting differentially methylated loci in three breast cancer cell lines. *Cancer Inform*, 3:43–54, 2007.
- [14] A. Khalili, D. Potter, P. Yan, L. Li, J. Gray, T. Huang, and S. Lin. Gamma-normal-gamma mixture model for detecting differentially methylated loci in three breast cancer cell lines. *Cancer Inform*, 3:43–54, 2007.
- [15] A. M. Khan, H. El-Daly, and N. M. Rajpoot. A gamma-gaussian mixture model for detection of mitotic cells in breast cancer histopathology images. In *ICPR*, pages 149–152. IEEE, 2012.

- [16] Z. Li, Q. Liu, J. Cheng, and H. Lu. A variational inference based approach for image segmentation. In *Pattern Recognition, 2008. ICPR 2008. 19th International Conference on*, pages 1–4, Dec 2008.
- [17] B. G. Lindsay. *Mixture Models: Theory, Geometry and Applications*. NSF-CBMS Conference series in Probability and Statistics, Penn. State University, 1995.
- [18] A. Llera and C. Beckmann. Gaussian / Inverse Gamma mixture models of ICA maps . In *21th International Conference on Functional Mapping of the Human Brain*, 2015.
- [19] A. Llera and C. F. Beckmann. Estimating an Inverse Gamma distribution. *ArXiv e-prints*, May 2016.
- [20] D. J. C. MacKay. *Information Theory, Inference and Learning Algorithms*. Cambridge University Press, 2003.
- [21] S. Makni, J. Idier, and J. B. Poline. Bayesian joint detection-estimation of brain activity using mcmc with a gamma-gaussian mixture prior model. In *Proc. 31th Proc. IEEE ICASSP*, 2006.
- [22] J. Marin, K. Mengersen, and C. Robert. *Bayesian Modelling and Inference on Mixtures of Distributions*, volume 25. Handbook of Statistics, Dey D. and Rao C.R., Elsevier Sciences, London, 2005.
- [23] R. Pruim, M. Mennes, D. van Rooij, A. Llera, J. Buitelaar, and C. Beckmann. Ica-aroma: A robust ica-based strategy for removing motion artifacts from fmri data. *Neuroimage*, 2015.
- [24] X. Shang. Radar detection based on compound-gaussian model with inverse gamma texture. *IET Radar, Sonar and Navigation*, 5:315–321(6), March 2011.
- [25] S. Smith, M. Jenkinson, M. Woolrich, C. Beckmann, T. Behrens, H. Johansen-Berg, P. Bannister, M. De Luca, I. Drobnjak, D. Flitney, R. Niazy, J. Saunders, J. Vickers, Y. Zhang, N. De Stefano, J. Brady, and P. Matthews. Advances in functional and structural mr image analysis and implementation as fsl. *Neuroimage*, 23 Suppl 1:S20819, 2004.
- [26] D. Von Rhein, H. Mennes, M. and van Ewijk, A. Groenman, M. Zwiers, J. Oosterlaan, D. Heslenfeld, B. Franke, P. Hoekstra, S. V. Faraone,



- C. Hartman, and J. Buitelaar. The neuroimage study: a prospective phenotypic, cognitive, genetic and mri study in children with attention-deficit/hyperactivity disorder. design and descriptives. *Eur. Child Adolesc. Psychiatry*, 24:265281, 2015.
- [27] J. Wang, Y. Xia, J. Wang, and D. D. Feng. Variational bayes inference based segmentation of heterogeneous lymphoma volumes in dual-modality PET-CT images. In *2011 International Conference on Digital Image Computing: Techniques and Applications (DICTA), Noosa, QLD, Australia, December 6-8, 2011*, pages 274–278, 2011.
- [28] M. Woolrich, T. Behrens, C. Beckmann, and S. Smith. Mixture models with adaptive spatial regularization for segmentation with an application to fmri data. *Medical Imaging, IEEE Transactions on*, 24(1):1–11, Jan 2005.
- [29] M. W. Woolrich and T. E. Behrens. Variational bayes inference of spatial mixture models for segmentation. *IEEE Transactions on Medical Imaging*, 25(10):1380–1391, Oct 2006.
- [30] M. W. Woolrich, S. Jbabdi, B. Patenaude, M. Chappell, S. Makni, T. Behrens, C. Beckmann, M. Jenkinson, and S. M. Smith. Bayesian analysis of neuroimaging data in {FSL}. *NeuroImage*, 45(1, Supplement 1):S173 – S186, 2009. Mathematics in Brain Imaging.
- [31] X. X. Shang and H. Song. Radar detection based on compound-gaussian model with inverse gamma texture. *IET Radar, Sonar and Navigation*, 5:315–321(6), March 2011.

## 5 Appendices

### A Variational mixture models

Here we continue with the notation and the problem described in sections 2.1 and 2.2. The joint probability density function is given by

$$\begin{aligned} p(\mathbf{x}, \mathbf{Z}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}) = \\ = p(\mathbf{x}|\mathbf{Z}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r})p(\mathbf{Z}|\boldsymbol{\pi})p(\boldsymbol{\pi})p(\mu_1)p(\tau_1) \prod_{k=2,3} p(s_k)p(r_k) \end{aligned} \quad (1)$$

The conditional distribution over  $\mathbf{Z}$  given the mixing coefficients  $\boldsymbol{\pi}$  is

$$p(\mathbf{Z}|\boldsymbol{\pi}) = \prod_{n=1}^N \prod_{k=1}^3 \pi_k^{z_{nk}}. \quad (2)$$

The conditional distribution of the observations given the latent variables and each component parameters is

$$\begin{aligned} p(\mathbf{x}|\mathbf{Z}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}) = \\ = \prod_{n=1}^N \left[ p_1(x_n|\mu_1, \tau_1)^{z_{n1}} \prod_{k=2}^3 p_2(x_n|s_k, r_k)^{z_{nk}} \right]. \end{aligned} \quad (3)$$

Now, we introduce the priors over the parameters  $\boldsymbol{\pi}, \mu_1, \tau_1, s_2, r_2, s_3, r_3$ . The prior over the mixing proportions is symmetric Dirichlet ( $\lambda_k = \lambda_0 \forall k \in \{1, 2, 3\}$ ),

$$p(\boldsymbol{\pi}) = \mathcal{D}(\boldsymbol{\pi}|\lambda_0) = C(\lambda_0) \prod_{k=1}^3 \pi_k^{\lambda_0-1}.$$

We use a Gaussian prior for the mean  $\mu_1$  of the Gaussian component, parametrized using mean  $m_0$  and precision  $\tau_0$ ,

$$p(\mu_1) = \mathcal{N}(\mu_1, |m_0, \tau_0),$$

and a Gamma prior, parametrized using shape  $c^0$  and scale  $b^0$ , for the precision  $\tau_1$

$$p(\tau_1) = \mathcal{G}_2(\tau_1 | c^0, b^0).$$

For the non-Gaussian components of the mixture model (second and third components) we use a Gamma prior over  $\mathbf{r} = (r_2, r_3)$ , parametrized using shape  $d_0$  and rate  $e_0$ ,

$$p(\mathbf{r}) = \prod_{k=2}^3 \mathcal{G}(r_k | d_0, e_0).$$

For the shape parameters  $\mathbf{s} = (s_2, s_3)$ , we use a prior of the form

$$p(\mathbf{s}) \propto \prod_{k=2}^3 p(s_k),$$

where

$$p(s_k) \propto \frac{a_0^{s_k-1} r_k^{s_k c_0}}{\Gamma(s_k)^{b_0}} \quad (4)$$

if component  $k$  is Gamma distributed<sup>2</sup> and

$$p(s_k) \propto \frac{a_0^{-s_k-1} r_k^{s_k c_0}}{\Gamma(s_k)^{b_0}} \quad (5)$$

if component  $k$  is inverse-Gamma distributed.

These functionals depend on the rates  $\mathbf{r}$  and on three hyper parameters  $(a_0, b_0, c_0)$ . Equation (4) is an unnormalized conjugate prior for the shape of the Gamma distribution [10] and Equation (5) an unnormalized conjugate prior for the shape parameter of an inverse-Gamma distribution [19].

## A.1 Variational updates

We consider a variational distribution that factorizes between latent variables and parameters as

$$q(\mathbf{Z}, \boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}) = q(\mathbf{Z})q(\boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}).$$

---

<sup>2</sup> $\Gamma$  denotes the Gamma function

### A.1.1 Latent variables

Given a data vector of observations  $\mathbf{x} = \{x_1, \dots, x_N\}$ ,  $x_i \in \mathbb{R}$  and using standard VB results, we have that

$$\log q^*(\mathbf{Z}) = \langle \log p(\mathbf{x}, \mathbf{Z}, \boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}) \rangle_{\boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}} + \text{const.}$$

Considering Equation (1) and keeping only terms that depend on  $\mathbf{Z}$ , we have that

$$\begin{aligned} \log q^*(\mathbf{Z}) &= \\ &= \langle \log p(\mathbf{Z}|\boldsymbol{\pi}) \rangle_{\boldsymbol{\pi}} + \langle [\log p(\mathbf{x}|\mathbf{Z}, \mu_1, \tau_1, \mathbf{r}, \mathbf{s})] \rangle_{\mu_1, \tau_1, \mathbf{r}, \mathbf{s}} + \text{const.} \end{aligned}$$

Substituting conditionals and absorbing terms that are independent from  $\mathbf{Z}$  into the constant term, we obtain

$$\log q^*(\mathbf{Z}) = \sum_{n=1}^N \sum_{k=1}^3 z_{nk} \log \rho_{nk} + \text{const}, \quad (6)$$

where

$$\begin{aligned} \log \rho_{n1} &= \langle \log \pi_1 \rangle + \frac{1}{2} \langle \log \tau_1 \rangle - \frac{1}{2} \log(2\pi) + \\ &\quad - \frac{1}{2} \langle (x_n - \mu_1)^2 \rangle_{\mu_1} \langle \tau_1 \rangle. \end{aligned}$$

For the Gamma components we have

$$\begin{aligned} \log \rho_{nk} &= \langle \log \pi_k \rangle + (\langle s_k \rangle - 1) \log(x_n) + \langle s_k \rangle \langle \log r_k \rangle + \\ &\quad - \langle \log \Gamma(s_k) \rangle - \langle r_k \rangle x_n. \end{aligned}$$

For inverse-Gamma components we have

$$\begin{aligned} \log \rho_{nk} &= \langle \log \pi_k \rangle - (\langle s_k \rangle + 1) \log(x_n) + \langle s_k \rangle \langle \log r_k \rangle + \\ &\quad - \langle \log \Gamma(s_k) \rangle - \frac{\langle r_k \rangle}{x_n}. \end{aligned}$$

Due to the positive support of the Gamma/inverse-Gamma distributions and the negative support of negative Gamma/inverse-Gamma distributions, we define  $\log \rho_{nk} = -\infty$  if  $x_n < 0$  and component  $k$  is positive or, if  $x_n > 0$  and component  $k$  is negative.

Exponentiating both sides of (6) we have

$$q^*(\mathbf{Z}) \propto \prod_{n=1}^N \prod_{k=1}^3 \rho_{nk}^{z_{nk}},$$

so

$$q^*(\mathbf{Z}) = \prod_{n=1}^N \prod_{k=1}^3 \gamma_{nk}^{z_{nk}},$$

where

$$\gamma_{nk} = \frac{\rho_{nk}}{\sum_{j=1}^3 \rho_{nj}}.$$

### A.1.2 Model parameters

Turning to the functional  $q(\boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r})$ , we now derive the VB updates for the parameters  $w \in \{\boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{r}, \mathbf{s}\}$ .

First, we define

$$N_k = \sum_{n=1}^N \gamma_{nk},$$

$$\bar{x}_k = \sum_{n=1}^N \gamma_{nk} x_n.$$

Taking the expectations over  $\mathbf{Z}$  we have that

$$\begin{aligned} \log q^*(\boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}) &= \langle \log p(\mathbf{x}, \boldsymbol{\theta}) \rangle_{\mathbf{Z}} + \text{const} = \\ &= \sum_{n=1}^N \left[ \langle z_{n1} \rangle \log p_1(x_n | \mu_1, \tau_1) + \sum_{k=2}^3 \langle z_{nk} \rangle \log p_k(x_n | s_k, r_k) \right] + \end{aligned}$$

$$\begin{aligned}
& + \log p(\boldsymbol{\pi}) + \langle \log p(\mathbf{Z}|\boldsymbol{\pi}) \rangle_{\mathbf{Z}} + \\
& + \log p(\mu_1) + \log p(\tau_1) + \sum_{k=2}^3 \log p(s_k, r_k). \tag{7}
\end{aligned}$$

This expression is used to derive the parameter updates in the following subsections. In particular, for a given parameter  $w \in \{\boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{r}, \mathbf{s}\}$ , we identify terms in (7) that depend on  $w$  to get an expression for  $\log q^*(w)$ . Exponentiating and regrouping terms lead us to the rest of the updates.

For  $\boldsymbol{\pi}$ , we have

$$q^*(\boldsymbol{\pi}) = \mathcal{D}(\boldsymbol{\pi}|\hat{\boldsymbol{\lambda}}),$$

$$\hat{\lambda}_k = \lambda_0 + N_k.$$

For  $\mu_1$ , we have

$$q^*(\mu_1) = \mathcal{N}(\mu|\hat{m}, \hat{\tau}),$$

$$\hat{m} = \frac{1}{\hat{\tau}} (\tau_0 m_0 + \langle \tau_1 \rangle \bar{\mathbf{x}}_1),$$

$$\hat{\tau} = \tau_0 + \langle \tau_1 \rangle N_1.$$

For  $\tau_1$ , we have

$$q^*(\tau_1) = \mathcal{G}(\tau_1|\hat{c}, \hat{b}),$$

$$\hat{b} = \left[ \frac{1}{b^0} + \frac{1}{2} \sum_{n=1}^N \gamma_{n1} (x_n^2 + \langle \mu_1^2 \rangle - 2x_n \langle \mu_1 \rangle) \right]^{-1},$$

$$\hat{c} = c^0 + \frac{1}{2} N_1.$$

For  $\mathbf{r}$ , we have

$$q^*(\mathbf{r}) = \prod_{k=2}^3 \mathcal{G}_2(r_k | \hat{d}_k, \hat{e}_k),$$

$$\hat{d}_k = d_0 + \langle s_k \rangle N_k,$$

$$\hat{e}_k = e_0 + \bar{\mathbf{x}}_k.$$

For  $\mathbf{s}$ , we have, for Gamma components

$$q^*(s_k) \propto \frac{\hat{a}_k^{s_k-1} r_k^{s_k \hat{c}_k}}{\Gamma(s_k)^{\hat{b}_k}},$$

and for inverse-Gamma components

$$q^*(s_k) \propto \frac{\hat{a}_k^{-s_k-1} r_k^{s_k \hat{c}_k}}{\Gamma(s_k)^{\hat{b}_k}}.$$

In both cases we have

$$\hat{a}_k = a_0 \prod_{n=1}^N x_n^{\gamma_{nk}},$$

$$\hat{b}_k = b_0 + N_k,$$

$$\hat{c}_k = c_0 + N_k.$$

## A.2 Computing the expectations

Using standard results for a Dirichlet distribution we have that for  $k \in \{1, 2, 3\}$  the posterior required expectations over  $\pi_k$  are given by

$$\langle \pi_k \rangle = \frac{\hat{\lambda}_k}{\sum_{m=1}^3 \hat{\lambda}_m},$$

$$\langle \log \pi_k \rangle = \Psi(\hat{\lambda}_k) - \Psi\left(\sum_{m=1}^3 \hat{\lambda}_m\right).$$

Using standard results for a Gaussian distribution we have that the required posterior expectations over  $\mu_1$  are

$$\langle \mu_1 \rangle = \hat{m},$$

$$\langle \mu_1^2 \rangle = \hat{m}^2 + \frac{1}{\hat{\tau}}.$$

Using standard results for a Gamma distribution we have that the required posterior expectations over  $\tau_1$  are

$$\langle \tau_1 \rangle = \hat{b}\hat{c},$$

$$\langle \tau_1^2 \rangle = \hat{b}\hat{c}(1 + \hat{c}),$$

$$\langle \log \tau_1 \rangle = \Psi(\hat{c}) + \log \hat{b},$$

and considering a Gamma distribution parametrized using shape and rate we obtain the required posterior expectations over  $r_k$  for  $k \in \{2, 3\}$

$$\langle r_k \rangle = \frac{\hat{d}_k}{\hat{e}_k},$$

$$\langle \log r_k \rangle = \Psi(\hat{d}) - \log \hat{e}.$$

We compute the required expectations over  $s$  using the Laplace approximation. Consider the prior on the Gamma shape with the form of equation (4),

$$p_G(s|a, b, c, r) \propto \frac{a^{s-1} r^{sc}}{\Gamma(s)^b}$$

and the prior on the inverse-Gamma shape with the form of equation (5),

$$p_{IG}(s|a, b, c, r) \propto \frac{a^{-s-1} r^{sc}}{\Gamma(s)^b}$$



Making use of the chain rule we have

$$\frac{d \log p(s|a, b, c, r)}{ds} = \frac{d \log p(s|a, b, c, r)}{dp(s|a, b, c, r)} \frac{dp(s|a, b, c, r)}{ds},$$

and, since,

$$\frac{dp_G(s|a, b, c, r)}{ds} = p_G(s)[\log a + c \log r - b\Psi(s)], \quad (8)$$

and

$$\frac{dp_{IG}(s|a, b, c, r)}{ds} = p_{IG}(s)[- \log a + c \log r - b\Psi(s)], \quad (9)$$

we have that

$$\frac{d \log p_G(s|a, b, c, r)}{ds} = \log a + c \log r - b\Psi(s)$$

and

$$\frac{d \log p_{IG}(s|a, b, c, r)}{ds} = - \log a + c \log r - b\Psi(s).$$

Further, both second derivatives are equal for both cases

$$\frac{d^2 \log p_G(s|a, b, c, r)}{d^2 s} = \frac{d^2 \log p_{IG}(s|a, b, c, r)}{d^2 s} = -b\Psi_1(s),$$

where  $\Psi_1(s) = \frac{d\Psi(s)}{ds}$ . Therefore

$$p_G(s|a, b, c, r) \approx \mathcal{N}(s|\mu_G, b\Psi_1(\mu))$$

and

$$p_{IG}(s|a, b, c, r) \approx \mathcal{N}(s|\mu_{IG}, b\Psi_1(\mu)),$$

where

$$\mu_G = \Psi^{-1}\left(\frac{\log a + c \log r}{b}\right)$$

is a zero of (8) and

$$\mu_{IG} = \Psi^{-1}\left(\frac{-\log a + c \log r}{b}\right)$$

is a zero of (9).

Using these approximations we have that the first required expectation is approximated in the case of the Gamma by

$$\langle s_k \rangle \approx \Psi^{-1}\left(\frac{\log a_k + c_k \log r_k}{b_k}\right)$$

and, in the one of the inverse-Gamma, by

$$\langle s_k \rangle \approx \Psi^{-1}\left(\frac{-\log a_k + c_k \log r_k}{b_k}\right).$$

The other required expectation is  $\mathbb{E}[\log(\Gamma(s_k))]$ . We use Taylor expansion to obtain

$$\mathbb{E}[\log(\Gamma(s))] \approx \mathbb{E}[\log(\Gamma(\mu))] + \frac{1}{b} + \frac{\Psi_2(\mu)\mu}{\Psi_1(\mu)b}.$$

### A.3 Hyper-parameters and initialization.

For the Gaussian component, we fixed the hyper-prior parameters values at  $m_0 = 0$ ,  $\tau_1 = 1$ ,  $c^0 = 0.01$  and  $b^0 = 100$ . This ensures that the mean is approximately centered at zero with a flat prior for the variance. For the Gamma (or inverse Gamma) components, we chose a prior distributions such that both mean and variance are set to 10. We then use the method of moments (see Appendix A) to estimate the prior distribution parameters ( $s_0$  shape and  $r_0$  rate/scale). We set  $d_0 = r_0$  and  $e_0 = 1$ , so that the variance on  $r_0$  has the same magnitude. For the hyper-priors on the shape parameter,  $s_0$ , we use the Laplace approximation (see Appendix A.2) to define a prior with the required expected value (and variance), resulting in

$$b_0 = c_0 = \frac{1}{s_0 \Psi_1(s_0)}.$$

For Gamma components, we have

$$\log a_0 = b_0 \Psi(s_0) - c_0 \log r_0.$$

For inverse-Gamma components, we have

$$\log a_0 = -b_0 \Psi(s_0) + c_0 \log r_0.$$

Finally the prior over the mixing proportions is fixed to  $\lambda_0 = 5$ .

The mixture model parameter initialization is performed using k-means [20]. The estimated means and variances are transformed into parameters for Gamma or inverse Gamma distributions for the required components using the method of moments (see Appendix B). These parameters are also used to estimate the density of each sample with respect to the non-Gaussian components required to estimate all initial  $\gamma_{nk}$ .

## A.4 Convergence

The convergence of the algorithms are monitored using the negative free energy (NFE). The NFE for the proposed model is given by

$$\begin{aligned} F = & \langle \log p(\mathbf{x}, \mathbf{Z} | \boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}) \rangle_{\mathbf{Z}, \boldsymbol{\pi}, \mu_1, \tau_1, \mathbf{s}, \mathbf{r}} + \\ & + \mathcal{H}[q^*(\mathbf{Z})] - \mathcal{KL}[\boldsymbol{\pi}] - \mathcal{KL}[\mu_1] - \mathcal{KL}[\tau_1] - \mathcal{KL}[\mathbf{s}] - \mathcal{KL}[\mathbf{r}]. \end{aligned}$$

The joint-likelihood (averaged over the posteriors) and the entropy term are straightforward to obtain. The  $\mathcal{KL}$ -divergences between priors and posteriors can be found elsewhere [7]. The only cumbersome term is  $\mathcal{KL}[\mathbf{s}]$  which does not have a known analytical solution. We therefore approximate the  $\mathcal{KL}$ -divergence by the  $\mathcal{KL}$ -divergence between the Gaussian approximations obtained by the Laplace approximations to  $p(\mathbf{s})$  (see Appendix A2).

## B Method of Moments

Given a data vector of observations  $\mathbf{x} = \{x_1, \dots, x_N\}$ ,  $x_i \in \mathbb{R}$ , and defining  $\mu$  as the mean of  $\mathbf{x}$  and  $v$  as its variance, the method of moments parameters estimation for the Gamma distribution reads

$$s \approx \frac{\mu^2}{v}, \quad \frac{1}{r} \approx \frac{v}{\mu}.$$

where  $s$  is the shape parameter and  $r$  is the rate parameter, and

$$s \approx \frac{\mu^2}{v} + 2, \quad r \approx \mu \left( \frac{\mu^2}{v} + 1 \right).$$

for the inverse-Gamma distribution where  $s$  is the shape parameter and  $r$  is the scale parameter.

## C State of the art mixture models

Algorithm 3 summarizes the approximated maximum likelihood algorithm presented in [5, 30] for learning a Gaussian/Gamma mixture model (GGM).

---

### Algorithm 3: ML Gauss Gamma mixture model (GGM)

---

**Require:** Data:  $\mathbf{x} = \{x_1, \dots, x_N\}, x_n \in \mathbb{R}$ ;

Parametrization:  $p(x_n|\Theta, \Pi) = \sum_{k=1}^K \pi_k p_k(x_n|\theta_k)$

$p_1(x_n|\Theta_1) = \mathcal{N}(x_n|\mu_1, v_1), p_2(x_n|\Theta_2) = IG(x_n|s_2, r_2), p_3(x_n|\Theta_K) = IG^-(x_n|s_3, r_3).$

1: Initialization parameter values:  $\Theta = \{\mu_1, v_1, s_2, r_2, s_3, r_3\}, \Pi = \{\pi_1, \pi_2, \pi_3\}$

2: **repeat**

3:     **for**  $n \in \{1, \dots, N\}$  **do**

4:         **for**  $k \in \{1, \dots, 3\}$  **do**

5:              $\gamma_k(x_n) = \frac{\pi_k p_k(x_n|\Theta_k)}{\sum_{j=1}^3 \pi_j p_j(x_n|\Theta_j)}.$

6:         **for**  $k \in \{1, \dots, 3\}$  **do**

7:              $\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_k(x_n) x_n$

8:              $v_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_k(x_n) (x_n - \mu_k)^2$

9:             **if**  $k \in \{2, 3\}$  **then**

10:                  $\alpha_k =$

11:                  $N_k = \sum_{n=1}^N \gamma_k(x_n)$

12:                  $\pi_k = \frac{N_k}{\sum_{j=1}^3 N_j}$

13: **until** convergence

14: **return**  $\Theta, \Pi.$

---

Algorithm 4 summarizes the algorithm presented in [18] for learning a Gaussian/inverse-Gamma mixture model (GIM).

---

**Algorithm 4: ML Gauss inverse-Gamma mixture model (GIM)**


---

**Require:** Data:  $\mathbf{x} = \{x_1, \dots, x_N\}, x_n \in \mathbb{R}$ ;

Parametrization:  $p(x_n|\Theta, \Pi) = \sum_{k=1}^K \pi_k p_k(x_n|\theta_k)$

$p_1(x_n|\Theta_1) = \mathcal{N}(x_n|\mu_1, v_1), p_2(x_n|\Theta_2) = IG(x_n|s_2, r_2), p_3(x_n|\Theta_K) = IG^-(x_n|s_3, r_3)$ .

1: Initialization parameter values:  $\Theta = \{\mu_1, v_1, s_2, r_2, s_3, r_3\}, \Pi = \{\pi_1, \pi_2, \pi_3\}$

2: **repeat**

3:     **for**  $n \in \{1, \dots, N\}$  **do**

4:         **for**  $k \in \{1, \dots, 3\}$  **do**

5:              $\gamma_k(x_n) = \frac{\pi_k p_k(x_n|\Theta_k)}{\sum_{j=1}^3 \pi_j p_j(x_n|\Theta_j)}$ .

6:         **for**  $k \in \{1, \dots, 3\}$  **do**

7:              $\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_k(x_n) x_n$

8:              $v_k = \frac{1}{N_k} \sum_{n=1}^N \gamma_k(x_n) (x_n - \mu_k)^2$

9:             **if**  $k \in \{2, 3\}$  **then**

10:                  $\alpha_k = \frac{\mu_k^2}{v_k} + 2, \beta_k = \mu_k (\frac{\mu_k^2}{v_k} + 1)$

11:              $N_k = \sum_{n=1}^N \gamma_k(x_n)$

12:              $\pi_k = \frac{N_k}{\sum_{j=1}^3 N_j}$

13: **until** convergence

14: **return**  $\Theta, \Pi$ .

---